

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2002-055693

(43)Date of publication of application : 20.02.2002

(51)Int.Cl.

G10L 13/06

(21)Application number : 2000-242068

(71)Applicant : SANYO ELECTRIC CO LTD

(22)Date of filing : 10.08.2000

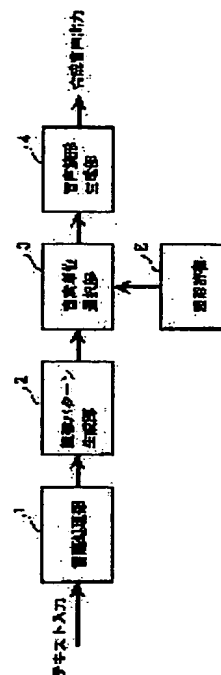
(72)Inventor : HIRAI HIROYUKI

(54) METHOD FOR SYNTHESIZING VOICE

(57)Abstract:

PROBLEM TO BE SOLVED: To provide a method for synthesizing a voice with which a phonemic piece with an inferior quality causing sound quality degradation is hardly selected as an optimum phonemic piece, without having to significantly correct waveform dictionary.

SOLUTION: The method for synthesizing voice consists of a step, in which penalty information is added to supplementary information for each phonemic unit, a step in which a user input a synthesized voice part of inferior quality when the user decides the quality of a synthesized voice as being inferior as a result of hearing, and a step in which the value for enlarging a calculated value distorted from a target forcibly, when the phonemic piece corresponding to the synthesized voice part of inferior quality is selected as a candidate is set to penalty information in the supplementary information of the above phonemic piece, when the synthesized voice part of inferior quality inputted by the user is inputted.



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2002-55693

(P2002-55693A)

(43) 公開日 平成14年2月20日 (2002.2.20)

(51) IntCl.⁷

G 1 0 L 13/06

識別記号

F I

G 1 0 L 5/04

キーワード* (参考)

F 5 D 0 4 5

審査請求 有 請求項の数 2 O L (全 6 頁)

(21) 出願番号 特願2000-242068 (P2000-242068)

(22) 出願日 平成12年8月10日 (2000.8.10)

(71) 出願人 000001889

三洋電機株式会社

大阪府守口市京阪本通2丁目5番5号

(72) 発明者 平井 啓之

大阪府守口市京阪本通2丁目5番5号 三

洋電機株式会社内

(74) 代理人 100086391

弁理士 香山 秀幸

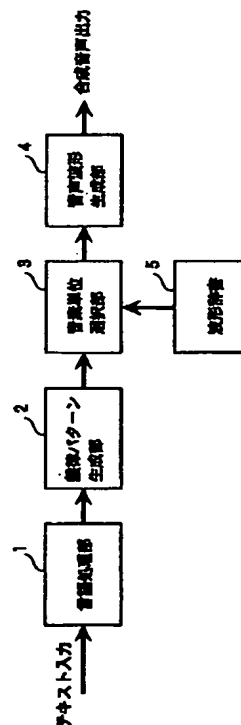
Fターム (参考) 5D045 AA20

(54) 【発明の名称】 音声合成方法

(57) 【要約】

【課題】 この発明は、波形辞書の大幅な修正を行なうことなく、音質劣化につながる品質の悪い音素片が最適な音素片として選択されにくくすることができる音声合成方法を提供することを目的とする。

【解決手段】 各音素単位の補助情報にペナルティ情報を追加しておくステップ、ユーザが音声合成結果を聞いて、その品質が悪い場合には、品質の悪い合成音声箇所をユーザに入力させるステップ、ならびにユーザによって入力された品質の悪い合成音声箇所に対応する音素片の補助情報内のペナルティ情報に、当該音素片が候補として選択されたときにターゲットとの歪み算出値を強制的に大きくさせるような値を設定するステップを備えている。



【特許請求の範囲】

【請求項1】 複数の音声単位と各音素単位毎にターゲットとの歪みを算出するために用いられる補助情報とが波形辞書に格納されており、波形辞書に格納されている音素単位の組み合わせの中で、ターゲットとの歪みが最も少なくなる組み合わせを選択する音素単位選択型の音声合成方法において、

各音素単位の補助情報にペナルティ情報を追加しておくステップ、

ユーザが音声合成結果を聞いて、その品質が悪い場合には、品質の悪い合成音声箇所をユーザに入力させるステップ、ならびにユーザによって入力された品質の悪い合成音声箇所が入力された場合には、当該品質の悪い合成音声箇所に対応する音素片の補助情報内のペナルティ情報に、当該音素片が候補として選択されたときにターゲットとの歪み算出値を強制的に大きくさせるような値を設定するステップ、

を備えていることを特徴とする音声合成方法。

【請求項2】 複数の音声単位と各音素単位毎にターゲットに対する適応度を算出するために用いられる補助情報とが波形辞書に格納されており、波形辞書に格納されている音素単位の組み合わせの中で、ターゲットに対する適応度が最も大きくなる組み合わせを選択する音素単位選択型の音声合成方法において、

各音素単位の補助情報に優先度情報を追加しておくステップ、ユーザが音声合成結果を聞いて、その品質が悪い場合には、品質の悪い合成音声箇所をユーザに入力させるステップ、ならびにユーザによって入力された品質の悪い合成音声箇所が入力された場合には、当該品質の悪い合成音声箇所に対応する音素片の補助情報内の優先度情報に、当該音素片が候補として選択されたときにターゲットに対する適応度の算出値を強制的に小さくさせるような値を設定するステップ、

を備えていることを特徴とする音声合成方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】この発明は、任意のテキスト情報を合成音声で読み上げることのできる音声合成方法に関する。

【0002】

【従来の技術】図1は、音声合成装置の概略構成を示している。

【0003】入力された日本語仮名漢字混じりのテキストは、言語処理部1で形態素解析、係り受け解析が行なわれ、音素記号、アクセント記号等に変換せしめられる。

【0004】韻律パターン生成部2では、音素記号、アクセント記号列および形態素解析結果から得られる入力テキストの品詞情報を用いて、音韻継続時間長（声の長

さ DUR^T ）、基本周波数（声の高さ FO^T ）、母音中心のパワー（声の大きさ POW^T ）等の推定が行なわれる。

【0005】音素単位選択部3では、推定された音韻継続時間長 DUR^T 、基本周波数 FO^T および母音中心のパワー POW^T に最も近く、かつ波形辞書5に蓄積されている音素単位（音素片）を接続したときの歪みが最も少なくなる音素片の組み合わせがDP（動的プログラミング）を用いて選択される。

【0006】音声波形生成部4では、選択された音素片の組み合わせにしたがって、ピッチを変換しつつ音素片の接続を行なうことによって音声が生産される。

【0007】図2は、波形辞書5の内容を示している。波形辞書5は、複数の音素片が格納された音素片格納部51と、音素片格納部51内の各音素片に関する補助情報が格納された補助情報格納部52とがある。補助情報には、音素片のパワー（ POW^{pic} ）、基本周波数（ FO^{pic} ）、継続時間長（ DUR^{pic} ）等がある。

【0008】ところで、音素単位選択部3では、波形辞書5に蓄積されている音素片の組み合わせの中で、歪みが少なくなる組み合わせを選択しているが、この歪みには次のようなものがある。

【0009】つまり、図3に示すように、 u_{i-1} 、 u_i 、 u_{i+1} を波形辞書5から抽出した音素片として、 t_{i-1} 、 t_i 、 t_{i+1} を実際に使用する環境（ターゲット）とすると、 u_i に対する歪みには、 C_i^t と、 C_i^c とがある。

【0010】ここで、 C_i^t は、 i 番目の音素について辞書から抽出した音素片（ u_i ）と実際に使用する環境（ターゲット t_i ）との間の歪みである。また、 C_i^c は、 i 番目の音素片（ u_i ）と、 $i-1$ 番目の音素片（ u_{i-1} ）とを接続したときに生じる歪みである。音素単位選択部3は、動的計画法（DP法）を用いて音素片を接続していき、入力された全ての音素に対する C_i^t と C_i^c との総和 C^{all} が最小となる音素片の組み合わせを選択する。

【0011】 C_i^t は、次の数式1で表される。

【0012】

【数1】

$$C_i^t = w_{pow}' D_{pow}'(t_i, u_i) + w_{fo}' D_{fo}'(t_i, u_i) + w_{dur}' D_{dur}'(t_i, u_i)$$

【0013】数式1において、各変数は、次のように定義される。

【0014】 $D_{pow}'(t_i, u_i)$ は、 i 番目の音素について、辞書から抽出した音素片（ u_i ）のパワー（ $POW^{pic}(i)$ ）と、実際に使用する環境（ターゲット t_i ）のパワー（ $POW^T(i)$ ）との間の距離の自乗であり、 $\{ (POW^{pic}(i)) - (POW^T(i)) \}^2$ となる。

【0015】 w_{pow}' は、 $D_{pow}'(t_i, u_i)$ に対する重み係数である。

【0016】 $D_{fo}'(t_i, u_i)$ は、 i 番目の音素について、辞書から抽出した音素片（ u_i ）の基本周波数

($FO^{Dic}(i)$)と、実際に使用する環境(ターゲット t_i)の基本周波数($FO^T(i)$)との間の距離の自乗であり、 $\{(FO^{Dic}(i)) - (FO^T(i))\}^2$ となる。

【0017】 w_{ro} は、 $D_{ro}(t_i, u_i)$ に対する重み係数である。

【0018】 $D_{ro}(t_i, u_i)$ は、 i 番目の音素について、辞書から抽出した音素片(u_i)の継続時間長($DUR^{Dic}(i)$)と、実際に使用する環境(ターゲット*

$$C_i = w_{row}^* D_{row}^*(u_i, u_{i-1}) + w_{ro}^* D_{ro}^*(u_i, u_{i-1}) + w_{src}^* D_{src}^*(u_i, u_{i-1})$$

【0022】数式2において、各変数は、次のように定義される。

【0023】 $D_{ro}(u_i, u_{i-1})$ は、 i 番目の音素片(u_i)の始端のパワー($POW^{Dic}(i)$)と、 $i-1$ 番目の音素片(u_{i-1})の終端のパワー($POW^{Dic}(i-1)$)との間の距離の自乗であり、 $\{(POW^{Dic}(i)) - (POW^{Dic}(i-1))\}^2$ となる。

【0024】 w_{ro} は、 $D_{ro}(u_i, u_{i-1})$ に対する重み係数である。

【0025】 $D_{ro}(u_i, u_{i-1})$ は、 i 番目の音素片(u_i)の始端の基本周波数($FO^{Dic}(i)$)と、 $i-1$ 番目の音素片(u_{i-1})の終端の基本周波数($FO^{Dic}(i-1)$)との間の距離の自乗であり、 $\{(FO^{Dic}(i)) - (FO^{Dic}(i-1))\}^2$ となる。

【0026】 w_{ro} は、 $D_{ro}(u_i, u_{i-1})$ に対する重み係数である。

【0027】 $D_{src}(u_i, u_{i-1})$ は、 i 番目の音素片(u_i)の始端のスペクトル($SPC^{Dic}(i, j)$, $j=1 \sim 16$)と、 $i-1$ 番目の音素片(u_{i-1})の終端のスペクトル($SPC^{Dic}(i-1, j)$, $j=1 \sim 16$)との間の距離の自乗であり、 $\{(SPC^{Dic}(i, j)) - (SPC^{Dic}(i-1, j))\}^2$ となる。

【0028】 w_{src} は、 $D_{src}(u_i, u_{i-1})$ に対する重み係数である。

【0029】入力された全ての音素に対する C_i と C_i との総和 C^{all} は、次の数式3で表される。

【0030】

【数3】

$$C^{all} = \sum_i C_i + C_i$$

【0031】

【発明が解決しようとする課題】ところで、上述したように音声合成方法によれば、品質の高い合成音声、つまり、自然発話に近い合成音声を得ることができる。しかしながら、自然発話から作成した音素片には、“なまけ”、“いいよども”など、実際に選択された場合に音質の劣化につながる音素片が存在している可能性が高い。このような音素片を含まないように波形辞書5を作成することが好ましいが、実際には音質劣化につながる音素片をすべて取り除いて波形辞書5を作成すること

*ト t_i)の継続時間長($DUR^T(i)$)との間の距離の自乗であり、 $\{(DUR^{Dic}(i)) - (DUR^T(i))\}^2$ となる。

【0019】 w_{dur} は、 $D_{dur}(t_i, u_i)$ に対する重み係数である。

【0020】 C_i は、次の数式2で表される。

【0021】

【数2】

は困難である。

【0032】また、波形辞書5を作成した後に、音質劣化につながる音素片を削除していくといったことも考えられるが、そのようにすると、波形辞書5の大幅な修正が必要となる。

【0033】この発明は、波形辞書の大幅な修正を行なうことなく、音質劣化につながる品質の悪い音素片が最適な音素片として選択されにくくすることができる音声合成方法を提供することを目的とする。

【0034】

【課題を解決するための手段】この発明による第1の音声合成方法は、複数の音声単位と各音素単位毎にターゲットとの歪みを算出するために用いられる補助情報とが波形辞書に格納されており、波形辞書に格納されている音素単位の組み合わせの中で、ターゲットとの歪みが最も少なくなる組み合わせを選択する音素単位選択型の音声合成方法において、各音素単位の補助情報にペナルティ情報を追加しておくステップ、ユーザが音声合成結果を聞いて、その品質が悪い場合には、品質の悪い合成音声箇所をユーザに入力させるステップ、ならびにユーザによって入力された品質の悪い合成音声箇所が入力された場合には、当該品質の悪い合成音声箇所に対応する音素片の補助情報内のペナルティ情報に、当該音素片が候補として選択されたときにターゲットとの歪み算出値を強制的に大きくさせるような値を設定するステップを備えていることを特徴とする。

【0035】この発明による第2の音声合成方法は、複数の音声単位と各音素単位毎にターゲットに対する適応度を算出するために用いられる補助情報とが波形辞書に格納されており、波形辞書に格納されている音素単位の組み合わせの中で、ターゲットに対する適応度が最も大きくなる組み合わせを選択する音素単位選択型の音声合成方法において、各音素単位の補助情報に優先度情報を追加しておくステップ、ユーザが音声合成結果を聞いて、その品質が悪い場合には、品質の悪い合成音声箇所をユーザに入力させるステップ、ならびにユーザによって入力された品質の悪い合成音声箇所が入力された場合には、当該品質の悪い合成音声箇所に対応する音素片の補助情報内の優先度情報に、当該音素片が候補として選択されたときにターゲットに対する適応度の算出値を強

制的に小さくさせるような値を設定するステップを備えていることを特徴とする。

【0036】

【発明の実施の形態】以下、この発明の実施の形態について説明する。

【0037】〔1〕第1の実施の形態の説明

音声合成装置の全体構成は、図1と同じである。

【0038】第1の実施の形態では、次の点(1)、(2)、(3)が従来と異なっている。

【0039】(1) 図4に示すように、各音素片の補*10

$$C_i' = w_{row}' D_{row}'(t_i, u_i) + w_{ro}' D_{ro}'(t_i, u_i) + w_{dur}' D_{dur}'(t_i, u_i) + D_{pi}'(u_i)$$

【0043】(3) ユーザが音声合成結果を聞いて、その品質が悪い場合には、品質の悪い合成音声箇所を音声合成装置に入力するようにする。音声合成装置は、ユーザによって入力された品質の悪い合成音声箇所が入力された場合には、品質の悪い合成音声箇所に対応する音素片の補助情報内のペナルティー情報 $D_{pi}'(u_i)$ の値を、所定値 α に設定する。

【0044】この所定値 α としては、たとえば、数式1の C_i' の予想される最大値の約100倍の値が用いられる。具体的には、任意数の文章を入力したときの数式1の最大値を実験により求めておき、その最大値の100倍の値を、所定値 α として設定する。

【0045】上記(1)、(2)、(3)のような変更を行なうことにより、ペナルティー情報 $D_{pi}'(u_i)$ の値として α が設定されている品質の悪い音素片(u_i)が候補として選択された場合には、その音素片とターゲットとの歪み C_i' が、従来法に比べて α 分だけ大きくなり、当該音素片(u_i)が最適な音素片として選択されにくくなる。

【0046】上記実施の形態によれば、波形辞書内に品質の悪い音素片が存在している場合に、その音素片を削除するといった大幅な辞書の修正を行なうことなく、音素片の補助情報にペナルティー情報 $D_{pi}'(u_i)$ を追加するといった小規模な修正を行なうことによって、品質の悪い音素片を選択されにくくすることができるようになる。

【0047】高品質の音声合成装置の場合には、波形辞書内の音素片格納部には6万個程度の音素片が格納されるため、音素片格納部の容量は数十MBに及ぶが、波形辞書内の補助情報格納部の容量は数MBというように、音素片格納部の容量の十分の1以下とである。このため、上記実施の形態のように補助情報格納部のみの修正を行なう方が容易である。また、音素片の削除に品質の改善を行なう従来方法では、波形辞書全てを置き換える必要があるが、上記実施の形態の方法では補助情報にペナルティー情報 $D_{pi}'(u_i)$ を追加するといった修※

$$S_i' = w_{row}'(1/D_{row}'(t_i, u_i)) + w_{ro}'(1/D_{ro}'(t_i, u_i)) + w_{dur}'(1/D_{dur}'(t_i, u_i))$$

【0055】また、数式5において、 S_i' は、i番目の音素片について辞書から選択した音素片(u_i)の始端

* 助情報に、ペナルティー情報 $D_{pi}'(u_i)$ を追加する。ペナルティー情報 $D_{pi}'(u_i)$ の初期値は、0である。

【0040】(2) 音素単位選択部3で歪み C_i' を算出するための C_i' に、ペナルティー情報 $D_{pi}'(u_i)$ をパラメータとして加える。

【0041】つまり、 C_i' は、次の数式4で表わされる。

【0042】

【数4】

※正のみであるため、波形辞書の一部の変更のみで修正が可能である。

【0048】また、ユーザが自由に波形辞書から品質の悪い音素片を削除することにより、合成音声の品質を改善させることも考えられるが、音素の種類によってはその音素に対応する全ての音素片を削除してしまうおそれがある。そうすると、当該音素を含む文章に対して合成音声を生成できなくなる可能性がある。

【0049】これに対して、上記実施の形態による方法では、たとえ、ある音素に対応する全ての音素片に対するペナルティー情報 $D_{pi}'(u_i)$ の値が所定値 α に設定されたとしても、当該音素を音声合成する際には、その音素に対応する音素片の中で最適な音素片が選択されるため、当該音素に対して合成音声を生成することができるという利点がある。

【0050】〔2〕第2の実施の形態の説明

第1の実施の形態においては、音素単位選択部3では、波形辞書に蓄積されている音素片の組み合わせの中で、歪みが少なくなる組み合わせを選択しているが、音素単位選択部として、波形辞書に蓄積されている音素片の組み合わせの中で、適応度が大きくなる組み合わせを選択するものが知られている。

【0051】適応度 S^{ad} は、一般的に次の数式5で表される。

【0052】

【数5】

$$S^{ad} = \sum_i S_i' + S_i'$$

【0053】数式5において S_i' は、i番目の音素について辞書から抽出した音素片(u_i)と実際に使用する環境(ターゲット t_i)との間の類似度を示しており、次の数式6で表される。数式6中の各変数は、数式1中の変数と同じである。

【0054】

【数6】

と、 $i-1$ 番目の音素について辞書から選択した音素片

(u_{i-1}) の終端との間の類似度を示しており、次の数

式7で表される。数式7中の各変数は、数式2中の変数*

$$S_i^c = w_{row}^c (1/D_{row}^c(u_i, u_{i-1})) + w_{ro}^c (1/D_{ro}^c(u_i, u_{i-1})) + w_{src}^c (1/D_{src}^c(u_i, u_{i-1}))$$

【0057】第2の実施の形態では、次の点(1)、

(2)、(3)が、適応度を用いて音素単位を選択する従来例と異なっている。

【0058】(1) 各音素片の補助情報に、優先度情報 $E_{pri}^c(u_i)$ を追加する。優先度情報 $E_{pri}^c(u_i)$ の初期値は、所定値である。

※10 【数8】

$$S_i^c = w_{row}^c (1/D_{row}^c(t_i, u_i)) + w_{ro}^c (1/D_{ro}^c(t_i, u_i)) + w_{dur}^c (1/D_{dur}^c(t_i, u_i)) + E_{pri}^c(u_i)$$

【0062】(3) ユーザが音声合成結果を聞いて、その品質が悪い場合には、品質の悪い合成音声箇所を音声合成装置に入力するようにする。音声合成装置は、ユーザによって入力された品質の悪い合成音声箇所が入力された場合には、品質の悪い合成音声箇所に対応する音素片の補助情報内の優先度情報 $E_{pri}^c(u_i)$ の値を、初期値より小さい値に設定する。

【0063】

【発明の効果】この発明によれば、波形辞書の大幅な修正を行なうことなく、音質劣化につながる品質の悪い音素片が最適な音素片として選択されにくくすることができ★

*と同じである。

【0056】

【数7】

※【0059】(2) 音素単位選択部3で適応度 S_i^{all}

を算出するための S_i^c に、優先度情報 $E_{pri}^c(u_i)$ をパラメータとして加える。

【0060】つまり、 S_i^c は、次式8で表わされる。

【0061】

【数8】

★きる。

【図面の簡単な説明】

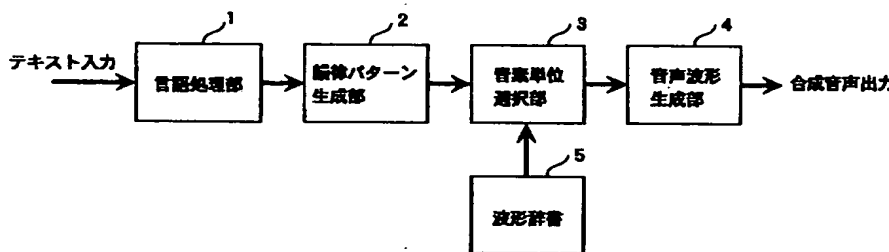
【図1】音声合成装置の全体構成を示すブロック図である。

【図2】波形辞書5の内容を示す模式図である。

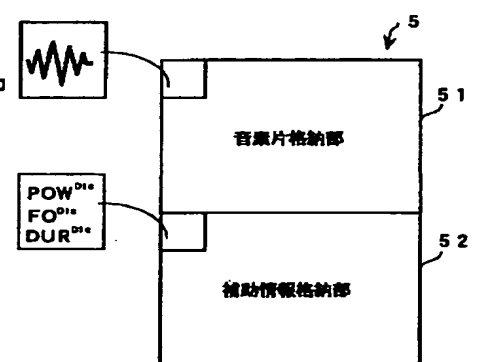
【図3】音素単位選択部3において、音素片の組み合わせを選択するために用いられる2種の歪み C_i^c 、 C_i^r を説明するための模式図である。

【図4】品質の悪い合成音声箇所に対応する音素片の補助情報に、ペナルティ情報 D_{pri}^c を追加された様子を示す模式図である。

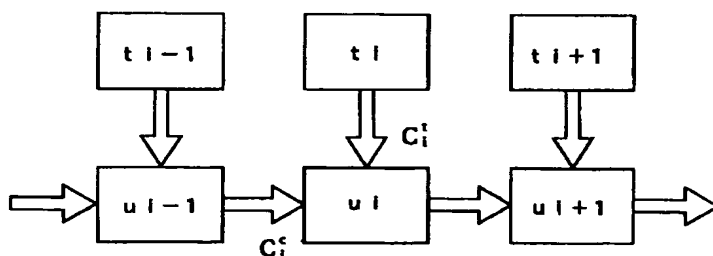
【図1】



【図2】



【図3】



【図4】

